



**UTM**  
UNIVERSITI TEKNOLOGI MALAYSIA

FACULTY OF COMPUTING  
SEMESTER II, SESSION 2024/2025

**BACHELOR OF COMPUTER SCIENCE (BIOINFORMATICS)**

**SECB3203 PROGRAMMING FOR BIOINFORMATICS - SECTION 01**

**PROJECT  
ALZHEIMER'S DISEASE PREDICTION**

GROUP MEMBERS	MATRIC NO
TAN ZHAO HONG	A23CS0188
CHIN PEI WEN	A23CS0065
KOO XUAN	A23CS0300

**LECTURER'S NAME : DR. SEAH CHOON SEN**

**SUBMISSION DATE : 03 DECEMBER 2025**

## **Table of Contents**

<b>3.0 Flowchart of the Proposed Approach.....</b>	<b>3</b>
<b>3.2 Model Development.....</b>	<b>3</b>
3.2.1 Imported Libraries.....	3
3.2.2 Data Preparation for Model Development.....	4
3.2.3 Classification models implemented for Alzheimer's disease prediction.....	5

## 3.0 Flowchart of the Proposed Approach

### 3.2 Model Development

For this project, classification was selected as the machine learning approach because the target variable, Alzheimer's disease diagnosis, is binary (0 = No Alzheimer's, 1 = Alzheimer's). Classification models are suitable for predicting discrete categories and are commonly used in medical diagnosis problems. Therefore, several classification algorithms were selected for model development and comparison.

#### 3.2.1 Imported Libraries

```
4  import numpy as np
5  import pandas as pd
6  import matplotlib.pyplot as plt
7  import seaborn as sns
8  from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
9  from sklearn.preprocessing import StandardScaler, LabelEncoder
10 from sklearn.metrics import (accuracy_score, precision_score, recall_score,
11                             f1_score, confusion_matrix, classification_report,
12                             roc_auc_score, roc_curve)
13 # from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.svm import SVC
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn.neighbors import KNeighborsClassifier
18 import warnings
19 warnings.filterwarnings('ignore')
```

Several Python libraries were imported to support data processing, model development, and visualization. Pandas and NumPy were used for data manipulation, Matplotlib and Seaborn for visualization, and Scikit-learn for model training, evaluation, and performance comparison.

### 3.2.2 Data Preparation for Model Development

```
255     # SPLIT AND SCALE
256     # split the data
257     X_train, X_test, y_train, y_test = train_test_split(
258         X, y, test_size=0.2, random_state=42, stratify=y
259     )
260
261     # Scale the features (Z-score normalization)
262     print("\n--- Final Scaling (Z-Score Normalization) ---")
263     scaler = StandardScaler()
264     X_train_scaled = scaler.fit_transform(X_train)
265     X_test_scaled = scaler.transform(X_test)
266
267     print(f"\nTraining set size: {X_train_scaled.shape}")
268     print(f"Test set size: {X_test_scaled.shape}")
269
270     print("\n" + "="*80)
271     print("DATA PREPROCESSING COMPLETE")
272     print("="*80)
273
274     return X_train_scaled, X_test_scaled, y_train, y_test, X.columns
275
```

The dataset was divided into independent variables (features) and the dependent variable (diagnosis). The data was then split into training and testing sets using an 80:20 ratio. Feature scaling was applied using StandardScaler to normalize the data and improve model performance.

### 3.2.3 Classification models implemented for Alzheimer's disease prediction.

```
276 # =====
277 # 4. MODEL TRAINING AND EVALUATION
278 # =====
279
280 def train_and_evaluate_models(X_train, X_test, y_train, y_test):
281     """Train multiple models and evaluate their performance"""
282
283     # Define models
284     models = {
285         'Logistic Regression': LogisticRegression(random_state=42, max_iter=1000),
286         'Decision Tree': DecisionTreeClassifier(random_state=42),
287         'SVM': SVC(random_state=42, probability=True),
288         'KNN': KNeighborsClassifier(n_neighbors=5)
289     }
290
291     results = {}
```

- **Logistic Regression**

Logistic Regression was implemented as a baseline classification model. It estimates the probability of a patient having Alzheimer's disease based on the input features. This model is widely used in binary classification problems due to its simplicity and interpretability.

- **Decision Tree**

Decision Tree was used to model decision rules in a tree-like structure. It is capable of capturing non-linear relationships and provides interpretability by visualizing decision paths.

- **Support Vector Machine (SVM)**

Support Vector Machine (SVM) was implemented to find the optimal hyperplane that separates Alzheimer's and non-Alzheimer's cases. SVM is effective in high-dimensional spaces and performs well in medical classification tasks.

- **K-Nearest Neighbors (KNN)**

K-Nearest Neighbors (KNN) classifies samples based on the majority class of the nearest neighbors. This model was included to compare instance-based learning with other classification techniques.

```

293     print("\n" + "="*80)
294     print("MODEL TRAINING AND EVALUATION")
295     print("*"*80)
296
297     for name, model in models.items():
298         print(f"\n{name}:")
299         print("-" * 40)
300
301         # Train the model
302         model.fit(x_train, y_train)
303
304         # Make predictions
305         y_pred = model.predict(x_test)
306         y_pred_proba = model.predict_proba(x_test)[:, 1] if hasattr(model, 'predict_proba') else None
307
308         # Calculate metrics
309         accuracy = accuracy_score(y_test, y_pred)
310         precision = precision_score(y_test, y_pred, average='weighted')
311         recall = recall_score(y_test, y_pred, average='weighted')
312         f1 = f1_score(y_test, y_pred, average='weighted')
313
314         # Store results
315         results[name] = {
316             'model': model,
317             'accuracy': accuracy,
318             'precision': precision,
319             'recall': recall,
320             'f1_score': f1,
321             'y_pred': y_pred,
322             'y_pred_proba': y_pred_proba
323         }
324
325         # Print results
326         print(f"Accuracy: {accuracy:.4f}")
327         print(f"Precision: {precision:.4f}")
328         print(f"Recall: {recall:.4f}")
329         print(f"F1-Score: {f1:.4f}")
330
331         # Cross-validation score
332         cv_scores = cross_val_score(model, x_train, y_train, cv=5, scoring='accuracy')
333         print(f"CV Accuracy: {cv_scores.mean():.4f} (+/- {cv_scores.std():.4f})")
334
335     return results

```

Figure : Model Training Loop

Each model is trained using the fit() method on the training set, then predictions are made on the test set for evaluation.